

Singapore Management University

Institutional Knowledge at Singapore Management University

Research Collection School Of Computing and Information Systems

School of Computing and Information Systems

4-2023

Open-set domain adaptation by deconfounding domain gaps

Xin ZHAO

Jilin University

Shengsheng WANG

Jilin University

Qianru SUN

Singapore Management University, qianrusun@smu.edu.sg

Follow this and additional works at: https://ink.library.smu.edu.sg/sis_research



Part of the [Artificial Intelligence and Robotics Commons](#), [Databases and Information Systems Commons](#), and the [Graphics and Human Computer Interfaces Commons](#)

Citation

ZHAO, Xin; WANG, Shengsheng; and SUN, Qianru. Open-set domain adaptation by deconfounding domain gaps. (2023). *Applied Intelligence*. 53, 7862-7875.

Available at: https://ink.library.smu.edu.sg/sis_research/7556

This Journal Article is brought to you for free and open access by the School of Computing and Information Systems at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection School Of Computing and Information Systems by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email cherylds@smu.edu.sg.

Open-Set Domain Adaptation by Deconfounding Domain Gaps

Xin Zhao^{a,b}, Shengsheng Wang^{a,b}, Qianru Sun^{c,*}

^a*College of Computer Science and Technology, Jilin University, Changchun 130012, China*

^b*Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, Jilin University, Changchun 130012, China*

^c*School of Computing and Information Systems, Singapore Management University*

Abstract

Open-Set Domain Adaptation (OSDA) aims to adapt the model trained on a source domain to the recognition tasks in a target domain while shielding any distractions caused by open-set classes, i.e., the classes “unknown” to the source model. Compared to standard DA, the key of OSDA lies in the separation between known and unknown classes. Existing OSDA methods often fail the separation because of overlooking the confounders (i.e., the domain gaps), which means their recognition of “unknown classes” is not because of class semantics but domain difference (e.g., styles and contexts). We address this issue by explicitly deconfounding domain gaps (DDP) during class separation and domain adaptation in OSDA. The mechanism of DDP is to transfer domain-related styles and contexts from the target domain to the source domain. It enables the model to recognize a class as known (or unknown) because of the class semantics rather than the confusion caused by spurious styles or contexts. In addition, we propose a module of ensembling multiple transformations (EMT) to produce calibrated recognition scores, i.e., reliable normality scores, for the samples in the target domain. Extensive experiments on two standard benchmarks verify that our proposed method outperforms a wide range of OSDA methods, because of its advanced ability of correctly recognizing unknown classes.

*Corresponding author

Email addresses: focusxin@outlook.com (Xin Zhao), wss@jlu.edu.cn (Shengsheng Wang), qianrusun@smu.edu.sg (Qianru Sun)

1. Introduction

Deep learning has made a remarkable success in a wide range of computer vision tasks [1, 2, 3], given a large amount of annotated training data. However, deep models can not generalize well to novel domains due to the domain shift [4]. To adapt these models, people always have to collect and annotate a large volume of training samples in the target domain as well, which is costly.

Unsupervised Domain adaptation (UDA) [5] tackles this issue by transferring knowledge from a source domain to a related but different domain (target domain) through using only unlabeled data. Most of UDA algorithms assume that the source and target datasets cover identical categories, known as Closed-Set Domain Adaptation (CSDA), as shown in Fig. 1a. While this assumption does not stand in real applications, as it is not possible to guarantee two domains sharing the same label space if no labels are available in one domain (the target domain). Therefore, researchers come up with a more reasonable and realistic setting called Open-Set Domain Adaptation (OSDA) [6, 7, 8]. The mainstream setting was introduced by Saito et al. [7], where the classes in the source domain are fully known and some of the classes in the target domain are unknown to the model trained in the source domain, as shown in Fig. 1b. The methods for OSDA specifically aim to classify the target domain samples correctly either into the label space of the source domain or as a special class called “unknown”.

The key in OSDA lies in how to effectively recognize and isolate the unknown samples, compared to the DA in closed-set scenarios. Existing methods usually define the normality score whose value shows to be lower for unknown sample than known sample. There are two typical issues. First, the model producing such scores is trained solely on source datasets, overlooking the confounder, i.e., the domain gap between source and target datasets. The essential reason behind is that when the model learns the semantic information of shared (known) classes, it is misled by spurious image styles or contexts. For example, if there

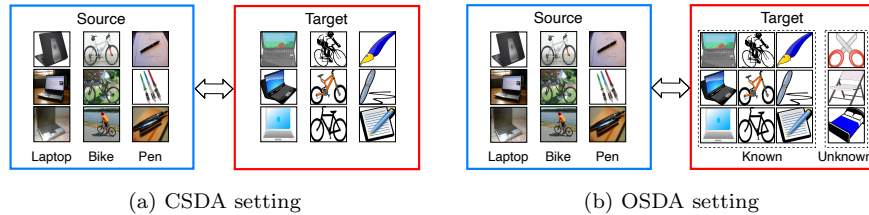


Figure 1: A comparison between the CSDA setting and the OSDA setting. (a) The CSDA setting assumes that the label space of two domains is identical. (b) The OSDA setting assumes that the target domain includes both known (shared) classes and unknown classes.

are many samples of “dog” on grass in source domain and “sheep” on grass in target domain (while there are quite fewer “dog” on grass in the target domain), the model gets misled by the context “on grass” and thus makes the wrong prediction on “sheep” samples as “dog” [9]. The second issue is that the model usually produces a single uncalibrated prediction on each input data, making the recognition of unknown samples unstable or unreliable.

In this paper, we solve the above issues in the two-stage framework presented in Fig. 2. In the first stage, we improve the ability and stability of the model to separate known and unknown samples. 1) We propose an explicit module of deconfounding domain gaps (DDP), which transfers image styles and contexts from the target domain to the source domain. We then fine-tune the model on the source samples with transferred styles and contexts to enable it to recognize samples as known (or unknown) because of their class semantics (rather than the confusion caused by spurious styles or contexts). 2) We propose a module of ensembling multiple transformations (EMT), which calibrates the model predictions by ensembling predictions from multiple transformations of each target sample. In the second stage, we leverage both the self-ensembling method [10] and the proposed DDP to deconfound domain gaps. Finally, we get the model that can recognize each target sample either as one of the known classes or as the special class “unknown”. We conduct extensive experiments on two OSDA benchmarks and show that both modules contribute to the performance improvement of the trained models. Therefore, our main contributions

are in three folds.

- (1) We point out that separating known and unknown classes remains a challenging problem in OSDA due to the confusion caused by domain gaps. We propose a novel OSDA method that can perform effective separation.
- 55 (2) We introduce an explicit module of deconfounding domain gaps (DDP), which transfers image styles and contexts from the target domain to the source domain and enables the model to correctly recognize unknown samples without confounded by domain gaps. In addition, we propose a module of ensembling multiple transformations (EMT) to calibrate the recognition
60 of the model and get more reliable normality scores.
- (3) We conduct experiments on two standard OSDA benchmarks. Our results demonstrate that our method outperforms the state-of-the-art. Our in-depth analyses verify that it gets better results because of its advanced ability of separating between known and unknown samples in the target
65 domain.

2. Related Works

In this section, we briefly review methods for domain adaptation and anomaly detection.

2.1. Domain Adaptation

70 **Closed-Set Domain Adaptation (CSDA)** works with the assumption that two domains have identical categories. CSDA approaches focus on mitigating the domain discrepancy between domains and can be grouped into several categories based on the adopted strategy. *Discrepancy-based* methods measure the divergence between domains in the feature space with a discrepancy metric,
75 such as Maximum Mean Discrepancy (MMD) [11, 12, 13], Higher-order moment matching (HoMM) [14], and Wasserstein distance [15]. The domain shift will be reduced by minimizing the metric during training. *Adversarial* methods

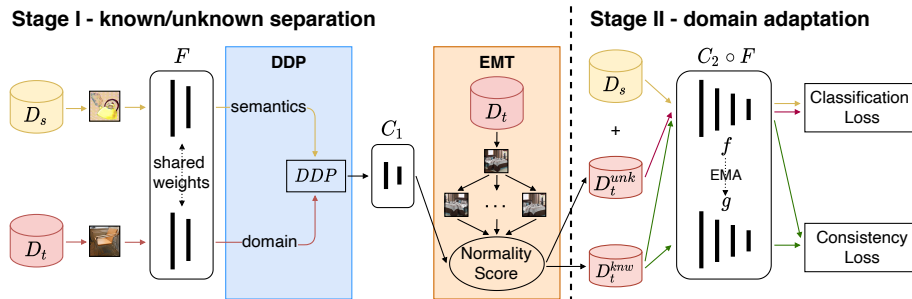


Figure 2: An illustration of the proposed method. Stage I: we propose an explicit module of deconfounding domain gaps (DDP), which transfers domain information (i.e., image styles and contexts) from the target domain to the source domain. We then train the encoder F and semantic network C_1 on the source samples with transferred styles and contexts to enable the model to recognize a class as known (or unknown) because of the class semantics rather than the confusion caused by spurious styles or contexts. After convergence, we propose a module of ensembling multiple transformations (EMT), which calibrates the model predictions by ensembling predictions from multiple transformations of each target sample. Based on the calibrated predictions, we can compute reliable normality scores used to divide the target datasets into a known target dataset \mathcal{D}_t^{knw} and an unknown target dataset \mathcal{D}_t^{unk} . Stage II: two networks with the same architecture are used: a student network $f = C_2 \circ F$, and a teacher network g with its weights being automatically set as an exponential moving average (EMA) of weights of the student network f . We minimize the consistency loss to mitigate the domain discrepancy between the source dataset \mathcal{D}_s and the target known dataset \mathcal{D}_t^{knw} . In addition, we can classify the target known samples and reject the target unknown samples by minimizing the classification loss on source dataset \mathcal{D}_s and unknown target dataset \mathcal{D}_t^{unk} . We also utilize the proposed DDP to further deconfound domain gaps, which is not shown in the right side of the figure for simplicity.

[16, 17, 18] leverage a domain classifier to distinguish source and target features, while training the feature encoder to device the domain classifier in order to extract domain-agnostic feature representations. Adversarial methods are the most popular ones and have obtained promising performance, which is further enhanced by recent works [19, 20, 21, 22, 23] with novel network designs. *Generative* methods [24, 25, 26, 27] leverage generative models to translate source samples to the target dataset and then reduce the domain discrepancy in both

85 feature and pixel levels. *Self-supervised* methods [28, 29, 30] design auxiliary self-supervised tasks for unlabeled target data to learn robust cross-domain representations. Consistency-enforcing methods [10, 31] force the model to make similar predictions for unannotated target samples even after they have been augmented.

90 **Open-Set Domain Adaptation (OSDA)** assumes target label set contains source label set. There are two different settings in the OSDA literature. Busto et al. [6] assumed each domain includes unknown categories besides the shared classes. And they proposed an algorithm called Assign-and-Transform-Iteratively (AIT), which maps target data to source domain and then utilizes
95 SVMs for final prediction. Saito et al. [7] eased the setting by requiring no unknown data from the source dataset, so target dataset contains all the source classes (known) and additional private classes that do not belong to the source (unknown). They also proposed a method, called Open Set Back-Propagation (OSBP), which adversarially trains a classifier with an extra 'unknown' class to
100 achieve common-private separation. Later OSDA methods all follow this more challenging and realistic setting. Separate To Adapt (STA) [8] aims to conduct known and unknown separation through a coarse-to-fine filtering process which includes two stages. First, multiple binary classifiers will be trained to compute the similarity score between source and target data. Second, target samples
105 with very low and high scores will be selected to train a final binary recognizer to distinguish known and unknown target samples. Attract or Distract (AoD) [32] leverage metric learning to match target samples with the corresponding neighborhood or distract away from the known classes. Rotation-based Open Set (ROS) [33] adopts rotation classification, a self-supervised method, to distinguish the known and unknown target samples and then adapt source knowledge
110 to the target known data.

Universal Domain Adaptation (UniDA), as a more general scenario, makes no assumption about the relationship of label sets between two domains. Universal Adaptation Network (UAN) [34] designs a measurement to evaluate
115 sample-level transferability based on domain similarity and prediction uncer-

tainty. Then samples with high transferability will be used with higher weight to promote common-class adaptation. However, as pointed by [35], this criterion is not discriminative and robust enough. Fu et al. [35] designed a better measurement that combines confidence, entropy, and consistency using multiple auxiliary classifiers to measure sample-wise uncertainty. Similarly, a class-level weighting strategy is applied for subsequent adversarial adaptation.

2.2. Anomaly Detection

Anomaly detection targets at detecting out-of-distribution (anomalous) samples by learning from normal samples. The approaches in this direction can be grouped into three categories. *Distribution-based* approaches [36, 37] leverage the normal samples to model the distribution function so that anomalous samples with lower likelihood can be filtered out. *Reconstruction-based* approaches [38, 39] leverage the encoder-decoder network to reconstruct the normal training samples. Then anomalous samples can be recognized as they have larger reconstruction error compared with normal samples. *Discriminative* approaches [40, 41] train a classifier on the normal samples and directly recognize anomalous samples based on the model prediction.

3. Method

In this section, we first formally introduce the preliminaries, then we present an overview of the proposed method and describe it in detail.

3.1. Preliminaries

We denote the annotated source domain drawn from distribution p_s as $\mathcal{D}_s = \{(x_j^s, y_j^s)\}_{j=1}^{N_s} \sim p_s$ and the unannotated target domain drawn from distribution p_t as $\mathcal{D}_t = \{x_j^t\}_{j=1}^{N_t} \sim p_t$. In OSDA, target label set \mathcal{C}_t contains source label set \mathcal{C}_s , i.e., $\mathcal{C}_s \subset \mathcal{C}_t$. We refer to classes from \mathcal{C}_s as the known classes and classes from $\mathcal{C}_t \setminus \mathcal{C}_s$ as the unknown classes. In OSDA, we both have $p_s \neq p_t$ and $p_s \neq p_t^{\mathcal{C}_s}$, where $p_t^{\mathcal{C}_s}$ represents the distribution of the target known data. Thus, we encounter both domain shift ($p_s \neq p_t^{\mathcal{C}_s}$) and class shift ($\mathcal{C}_s \neq \mathcal{C}_t$) problems

in OSDA. The goal of OSDA methods is to classify target known data correctly
145 and reject target unknown data.

OSDA introduces two challenges: negative transfer and known/unknown separation. (1) Enforcing to match the whole distribution of two domains as done in closed-set scenario will incur negative transfer, as the unknown target samples will also align mistakenly with source data. To solve this problem, we
150 need to apply adaptation only to the shared \mathcal{C}_s categories, mitigating the domain shift between p_s and $p_t^{\mathcal{C}_s}$. (2) Thus, we encounter the second challenge: known/unknown separation. All target samples should be recognized from target private categories $\mathcal{C}_t \setminus \mathcal{C}_s$ (unknown) or the shared categories \mathcal{C}_s (known).

3.2. Overview

155 To handle the aforementioned two challenges, we propose a novel OSDA method with a two-stage structure (Fig. 2): (i) we divide target datasets into known and unknown; (ii) we apply adaptation to source samples and target samples predicted as known. If we consider the unknown samples as anomalies, the first stage can be seen as an anomaly detection issue. And we can
160 also treat the second stage as a CSDA issue between target known and source distributions. Specifically, in the first stage, we propose an explicit module of deconfounding domain gaps (DDP), which transfers image styles and contexts from the target domain to the source domain, eliminating the confounding effect caused by domain gaps. In addition, we propose a module of ensembling multiple transformations (EMT) to calibrate the model predictions. Thus, we can
165 obtain more reliable normality scores based on the calibrated predictions. In the second stage, on the one hand, we leverage both the self-ensembling method [10] and the proposed DDP to reduce the domain discrepancy between the source data and target known data. On the other hand, we train the network to classify target known samples and reject target unknown samples by minimizing
170 the classification loss on source data and unknown target data.

3.3. Deconfounding Domain Gaps (DDP)

Recent domain generalization works observe that image styles and contexts are closely related to visual domains [42, 43]. Inspired by this observation, we propose an explicit module of deconfounding domain gaps (DDP) that transfers image styles and contexts from the target domain to the source domain. It enables the model to recognize a class as known (or unknown) because of the class semantics rather than the confusion caused by spurious styles or contexts. Following the common practice [44, 45], we use the feature statistics that preserved at the lower layers of the CNN as domain-related representation (i.e., styles and contexts) and their spatial configuration as semantic representation.

For an input sample x , we first obtain its feature maps $z \in \mathbb{R}^{C \times H \times W}$ from the feature encoder, where C indicates the number of channels, H and W represent spatial dimensions. Then we compute the channel-wise mean and standard deviation $\mu(z), \sigma(z) \in \mathbb{R}^C$ as style/context representation:

$$\mu(z) = \frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W z_{hw}, \quad (1)$$

$$\sigma(z) = \sqrt{\frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W (z_{hw} - \mu(z))^2 + \epsilon}. \quad (2)$$

Given intermediate feature maps $z^s, z^t \in \mathbb{R}^{C \times H \times W}$ corresponding to a source sample x^s and a target sample x^t , we replace the style/context of x^s with the style/context of x^t through adaptive instance normalization (AdaIN) [44]:

$$\text{DDP}(z^s, z^t) = \sigma(z^t) \cdot \left(\frac{z^s - \mu(z^s)}{\sigma(z^s)} \right) + \mu(z^t). \quad (3)$$

After training, the learned model can reject unknown samples only based on the semantic content without the influence of the confounder (i.e., the domain gaps).

3.4. Ensembling Multiple Transformations (EMT)

Previous methods usually utilize confidence (i.e., the largest probability of all classes) [7, 8] and uncertainty (i.e., entropy) [34, 33] as the normality score to

195 separate known (normal) and unknown (anomalous) samples, with an assumption that known samples have high confidence and low uncertainty, and vice versa. However, the confidence and the uncertainty used before are based on uncalibrated prediction, meaning they cannot represent the real confidence and uncertainty of the sample. Therefore, all the previous normality scores are not reliable and thus unable to separate known and unknown samples accurately.

To obtain more reliable normality score for separation, we propose a module of **ensembling multiple transformations (EMT)**, which ensembles predictions from multiple transformations of each target sample to calibrate the confidence and the entropy. Specifically, given a target image x^t , we apply random transformations (i.e., random crop and horizontal flip) to it to obtain m augmented samples $\{\tilde{x}_i^t\}_{i=1}^m$. $\hat{y}_i^t = C_1(F(\tilde{x}_i^t))$, ($i = 1, \dots, m$) is the corresponding prediction of each augmented sample \tilde{x}_i^t , where F and C_1 are the feature encoder and the semantic network. We compute the confidence w_{conf} and the entropy w_{ent} as follows:

$$w_{\text{conf}}(\hat{y}_i^t|_{i=1}^m) = \frac{1}{m} \sum_{i=1}^m \max(\hat{y}_i^t), \quad (4)$$

$$w_{\text{ent}}(\hat{y}_i^t|_{i=1}^m) = \frac{1}{m} \sum_{i=1}^m \left(\sum_{k=1}^{|\mathcal{C}^s|} -\hat{y}_{ik}^t \log(\hat{y}_{ik}^t) \right), \quad (5)$$

where \hat{y}_{ik}^t indicates the probability of k -th class and \max get the maximum entry in \hat{y}_i^t . We unify the w_{conf} and w_{ent} within $[0, 1]$ by the minmax normalization. The formulation of the normality score is:

$$\mathcal{N}(x^t) = \max\{w_{\text{conf}}, 1 - w_{\text{ent}}\}. \quad (6)$$

We maximize over these two terms to obtain the most reliable measurement.

3.5. Training Procedure

Stage I: known/unknown separation. To separate the known and unknown samples of \mathcal{D}_t , a CNN is trained on the source samples with transferred styles and contexts. To boost the discriminability of the model and facilitate

the following known/unknown separation, we also exploit the label smoothing (LS) as it pushes samples to distribute in tight evenly separated clusters [46]. The network consists of a feature encoder F and a semantic network C_1 . We train network by minimizing the following cross-entropy objective:

$$L_{\text{cls}} = -\mathbb{E}_{(x^s, y^s) \in \mathcal{D}_s, x^t \in \mathcal{D}_t} [y^{ls} \log C_1 (\text{DDP} (F(x^s), F(x^t)))] , \quad (7)$$

where $y^{ls} = (1 - \alpha)y^s + \alpha/|\mathcal{C}_s|$ indicates the smoothed label and α represents the smoothing parameter. After training, we compute the normality score for each target sample using F and C_1 as Eq. 6. Known samples have large values of \mathcal{N} , and vice versa. The target dataset can be divided into an unknown target dataset \mathcal{D}_t^{unk} and a known target dataset \mathcal{D}_t^{knw} using the normality score. We use the average of the normality score over all target samples $\bar{\mathcal{N}} = \frac{1}{N_t} \sum_{j=1}^{N_t} \mathcal{N}_j$ as the threshold, without the need to introduce any further parameter:

$$\begin{cases} x^t \in \mathcal{D}_t^{knw} & \text{if } \mathcal{N}(x^t) > \bar{\mathcal{N}} \\ x^t \in \mathcal{D}_t^{unk} & \text{if } \mathcal{N}(x^t) < \bar{\mathcal{N}}. \end{cases} \quad (8)$$

200 The detailed process about the computation of \mathcal{N} and the generation of \mathcal{D}_t^{knw} and \mathcal{D}_t^{unk} is described in Algorithm 1.

Stage II: domain adaptation. The problem is simplified to a CSDA problem after the target unknown data have been filtered out. Without the distraction of \mathcal{D}_t^{unk} , we can exploit \mathcal{D}_t^{knw} to decrease the domain discrepancy directly. In addition, \mathcal{D}_t^{unk} can be used to train the classifier to recognize the unknown samples. The network has a similar architecture to that of Stage I, consisting of a feature extractor F and a semantic network C_2 . The semantic network C_2 is the same as C_1 except for the last layer: the output dimension of C_1 is $|\mathcal{C}_s|$, while the output dimension of C_2 is $(|\mathcal{C}_s| + 1)$ because of the additional unknown class. We utilize the self-ensembling method [10] to close the domain gap. Two networks with the same architecture are used: a student network $f(x) = C_2(F(x))$, and a teacher network $g(x)$ with its weights being automatically set as an exponential moving average (EMA) of weights of the student network. The student network is trained to minimize the classification

Algorithm 1 Compute normality score and Generate \mathcal{D}_t^{knw} & \mathcal{D}_t^{unk}

Input:

Trained networks F and C_1

Target dataset $\mathcal{D}_t = \{x_j^t\}_{j=1}^{N_t}$

Output:

Known target dataset $\mathcal{D}_t^{knw} = \{x_j^{t, knw}\}_{j=1}^{N_t, knw}$

Unknown target dataset $\mathcal{D}_t^{unk} = \{x_j^{t, unk}\}_{j=1}^{N_t, unk}$

procedure GETENTROPYSCORE(y)

return $\sum_{k=1}^{|C^s|} -y_k \log(y_k)$

procedure GETNORMALITYSCORE(F, C_1, \mathcal{D}_t)

for each x_j^t **in** \mathcal{D}_t **do**

 Initialize: $\text{conf} = \{\}, \text{ent} = \{\}$

for each i **in** $\{1, \dots, m\}$ **do**

$\tilde{x}_j^t = \text{Transform}(x_j^t)$ # Apply transformation to x_j^t

$\hat{y}_j^t = C_1(F(\tilde{x}_j^t))$

$\text{conf} \leftarrow \max(\hat{y}_j^t)$

$\text{ent} \leftarrow \text{getEntropyScore}(\hat{y}_j^t)$

$w_{\text{conf}} = \text{mean}(\text{conf})$

$w_{\text{ent}} = \text{mean}(\text{ent})$

$w_{\text{conf}} = \text{normalize}(w_{\text{conf}})$ # Apply the minmax normalization

$w_{\text{ent}} = \text{normalize}(w_{\text{ent}})$

$\mathcal{N} \leftarrow \eta_j = \max\{w_{\text{conf}}, 1 - w_{\text{ent}}\}$

return \mathcal{N}

procedure MAIN()

 Initialize: $\mathcal{D}_t^{knw} = \{\}, \mathcal{D}_t^{unk} = \{\}$

$\mathcal{N} = \text{getNormalityScore}(F, C_1, \mathcal{D}_t)$

for each (x_j, η_j) **in** $(\mathcal{D}_t, \mathcal{N})$ **do**

if $\eta_j \geq \text{mean}(\mathcal{N})$ **then**

$\mathcal{D}_t^{knw} \leftarrow \mathbf{x}_j$

else

$\mathcal{D}_t^{unk} \leftarrow \mathbf{x}_j$

loss on source and target unknown samples, while maintaining consistent predictions with the teacher network for target known samples. The loss function of consistency can be formulated as:

$$L_{\text{con}} = \mathbb{E}_{x^t \in \mathcal{D}_t^{\text{knw}}} \left[(f(x^t) - g(x^t))^2 \right]. \quad (9)$$

The classification losses for samples from source and target unknown datasets are:

$$L_{\text{cls}}^s = -\mathbb{E}_{(x^s, y^s) \in \mathcal{D}_s, x^t \in \mathcal{D}_t^{\text{knw}}} \left[y^s \log C_2(\text{DDP}(F(x^s), F(x^t))) \right], \quad (10)$$

$$L_{\text{cls}}^{\text{unk}} = -\mathbb{E}_{(x^t, y^t) \in \mathcal{D}_t^{\text{unk}}} \left[y^t \log f(x^t) \right]. \quad (11)$$

It is worth noting that we also exploit the proposed DDP to transfer styles and contexts from the known target datasets to the source datasets, aiming to further deconfound domain gaps. We train the network to minimize the following overall objective:

$$L = (L_{\text{cls}}^s + L_{\text{cls}}^{\text{unk}}) + \lambda L_{\text{con}}, \quad (12)$$

where λ is the weight that trades off between classification loss and consistency loss. Once the training is complete, we predict the labels for all target samples using F and C_2 .

205 4. Experiments

In this section, we first introduce the experimental settings including datasets, compared approaches, evaluation metrics, and implementation details. Then, we present classification results on two standard datasets. Finally, we conduct further analysis to verify the effectiveness of the proposed method.

210 4.1. Experimental Settings

Datasets. **Office-31** [47] contains images within 31 classes collected from three visually different domains: *Webcam* (**W**) with 795 low-quality images obtained by web camera, *DSLR* (**D**) with 534 high-quality images taken by digital

SLR camera, and *Amazon* (**A**) with 2820 images obtained from *amazon.com*.
 215 Following the protocol in [7], we set the first 10 categories (1-10) as known and
 the last 11 categories (21-31) as unknown (in alphabetic order). We show some
 example images from Office-31 dataset in Fig. 3a. **Office-Home** [48] contains
 15,500 images within 65 classes collected from four different domains, Artistic
 images (**Ar**), Product images (**Pr**), Clip-Art images (**Cl**), and Real-World im-
 220 ages (**Rw**). Following the protocol in [8], we set the first 25 categories (1-25)
 in alphabetical order as known classes and the remaining 40 categories (26-65)
 as unknown. Office-Home is much more challenging than Office-31 due to the
 numerous categories and the large domain discrepancy. Some example images
 from Office-Home dataset are shown in Fig. 3b.

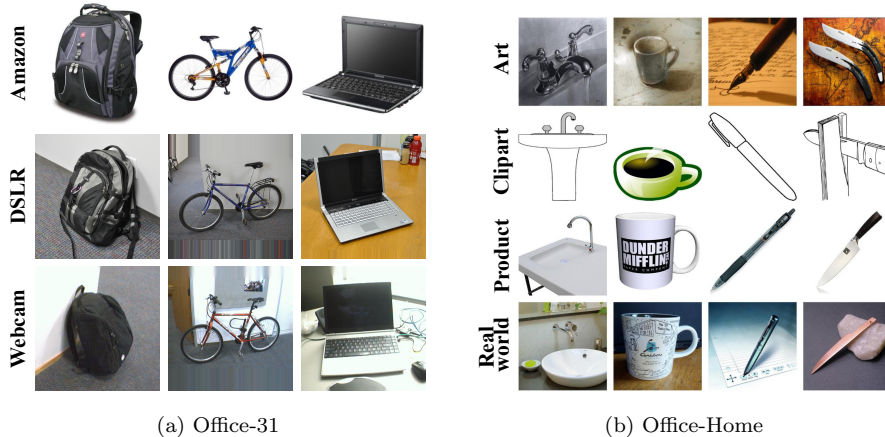


Figure 3: Example images in Office-31 and Office-Home.

225 **Compared Approaches.** We compare our method with: (1) Source Only
 model: ResNet-50 [1]; (2) CSDA method: DANN [16]; (3) OSDA methods:
 STA[8], OSBP[7], and ROS[33]; (4) UniDA method: UAN[34]. For ResNet-
 50 and DANN, we leverage a confidence threshold to separate unknown and
 unknown samples. All the results reported are the average over three random
 230 runs.

Evaluation Metrics. OS^* and UNK are two usual metrics used to eval-
 uate OSDA. OS^* denotes the average accuracy on known classes, and UNK

denotes the accuracy on the unknown class. They can be combined in $OS = \frac{|C_s|}{|C_s|+1} \times OS^* + \frac{1}{|C_s|+1} \times UNK$ to evaluate the overall performance. However, OS
 235 is not an appropriate metric as it assumes the accuracy of each known class has the same importance as the whole "unknown" class. Considering the trade-off between the accuracy of known and unknown classes is important in evaluating OSDAs methods, we exploit a metric: $HOS = 2 \frac{OS^* \times UNK}{OS^* + UNK}$ [33], which is the harmonic mean of OS^* and UNK . Unlike OS, HOS gives a high score only if
 240 the method achieves high performance both for known and unknown data.

Implementation Details. We utilize ResNet-50 [1] pretrained on ImageNet [49] as the backbone network. The feature encoder F consists of the first two residual blocks, while the remaining part combines the semantic network C_1 . We use the same hyperparameters for each dataset. Following DANN
 245 [16], we adjust the learning rate with $lr_p = \frac{lr_0}{(1+\omega p)^\phi}$, where p changes from 0 to 1 during the training process, lr_0 equal to 0.01 and 0.003 for Stage I and Stage II respectively, $\omega = 10$, and $\phi = 0.75$. The batch size is 32 for both two stages. For all the pretrained layers, the learning rate is 10 times lower than the layers learned from scratch. We adopt SGD to optimize the network, setting
 250 the momentum as 0.9 and the weight decay as 0.0005. In Stage I, we set the smoothing parameter to 0.1 and the number of multi-transformations (m) to 5. In Stage II, the trade-off parameter for consistency loss is $\lambda = 3$. We use the network learned in Stage I as the start for Stage II. The learning rate of the new unknown class is set to two times of the known classes.

255 4.2. Classification Results

To evaluate the performance of the OSDAs methods, we focus on the HOS as it can balance the importance between the accuracy of known (OS^*) and unknown classes (UNK), as discussed in Section 4.1. For a fair comparison, all results of the compared methods are either taken from [33] or obtained by
 260 running the code of [50].

Table 1 reports the classification results of Office-31. Our method outperforms all comparison approaches on most tasks except $W \rightarrow D$. Specifically,

Table 1: Accuracy (%) of all methods on Office-31 dataset.

	Office-31																				
	A → W			A → D			D → W			W → D			D → A			W → A			Avg.		
	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS
ResNet	67.7	65.9	66.8	78.5	62.8	69.8	93.6	73.0	82.0	98.6	79.3	87.9	58.1	81.0	67.7	56.9	80.6	66.7	75.6	73.8	73.5±1.2
DANN	93.8	62.9	75.3	88.5	63.3	73.8	98.1	42.7	59.5	99.8	47.9	64.7	68.4	51.5	58.7	71.9	56.7	63.4	86.8	54.2	65.9±1.4
STA	86.7	67.6	75.9	91.0	63.9	75.0	94.1	55.5	69.8	84.9	67.8	75.2	83.1	65.9	73.2	66.2	68.0	66.1	84.3	64.8	72.5±0.8
OSBP	86.8	79.2	82.7	90.5	75.5	82.4	97.7	96.7	97.2	99.1	84.2	91.1	76.1	72.3	75.1	73.0	74.4	73.7	87.2	80.4	83.7±0.4
UAN	95.5	31.0	46.8	95.6	24.4	38.9	99.8	52.5	68.8	81.5	41.4	53.0	93.5	53.4	68.0	94.1	38.8	54.9	93.4	40.3	55.1±1.4
ROS	88.4	76.7	82.1	87.5	77.8	82.4	99.3	93.0	96.0	100.0	99.4	99.7	74.8	81.2	77.9	69.7	86.6	77.2	86.6	85.8	85.9±0.2
Ours	89.5	79.0	83.9	86.4	82.7	84.5	99.8	96.3	98.0	100.0	98.8	99.4	75.2	84.3	79.5	70.1	90.2	78.9	86.8	88.6	87.4±0.4

our method significantly outperforms OSBP by 3.7%. Our method also boosts the HOS of state-of-the-art method ROS by 1.5%. In addition, we observe that
265 DANN, STA, and UAN perform even worse than the ResNet backbone since they suffer from negative transfer caused by the mismatching between the source samples and target unknown samples. The failure of these methods mainly due to their poor ability for target known and unknown separation.

We also compare our method with previous works on the challenging Office-
270 Home dataset. From Table 2, we can find that our method outperforms all compared methods on a total of 9 out of 12 transfer scenarios, demonstrating that our method works well with large domain gaps. On average, our method achieves the highest performance, 1.8% higher than the second-best method ROS. In addition, Our method outperforms STA and OSBP by a large margin,
275 6.9% and 3.3% respectively. The encouraging results indicate that our method is very effective for the OSDA setting.

From Tables 1 and 2, we can get one key observation that the advantage of our method is mainly due to its capability in distinguishing known and unknown samples. We can observe that while the average OS* of the compared methods is
280 close to ours, the UNK of our method is much higher, e.g., 2.8% and 3.4% higher than ROS on Office-31 and Office-Home respectively. This observation proves that our method is very significant for separating target known and unknown samples.

Table 2: Accuracy (%) of all methods on Office-Home dataset.

	Office-Home																										
	Pr → Rw			Pr → Cl			Pr → Ar			Ar → Pr			Ar → Rw			Ar → Cl											
	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS			
<i>ResNet</i>	70.2	61.0	65.3	32.8	67.1	44.1	45.7	70.5	55.5	64.4	64.0	64.2	76.9	59.7	67.2	44.8	68.7	54.2									
<i>DANN</i>	77.4	48.4	59.5	50.5	49.9	50.2	61.6	54.5	57.8	71.0	38.9	50.2	75.0	50.1	60.1	54.6	48.8	51.6									
<i>STA</i>	76.2	64.3	69.5	44.2	67.1	53.2	54.2	72.4	61.9	68.0	48.4	54.0	78.6	60.4	68.3	46.0	72.3	55.8									
<i>OSBP</i>	76.2	71.7	73.9	44.5	66.3	53.2	59.1	68.1	63.2	71.8	59.8	65.2	79.3	67.5	72.9	50.2	61.1	55.1									
<i>UAN</i>	84.0	0.1	0.2	59.1	0.0	0.0	73.7	0.0	0.0	81.1	0.0	0.0	88.2	0.1	0.2	62.4	0.0	0.0									
<i>ROS</i>	70.8	78.4	74.4	46.5	71.2	56.3	57.3	64.3	60.6	68.4	70.3	69.3	75.8	77.2	76.5	50.6	74.1	60.1									
<i>Ours</i>	69.3	76.9	72.9	48.6	75.6	59.2	56.3	68.3	61.7	65.5	79.4	71.8	76.4	78.2	77.3	50.1	83.9	62.7									

	Rw → Ar			Rw → Pr			Rw → Cl			Cl → Rw			Cl → Ar			Cl → Pr			Avg.		
	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS			
<i>ResNet</i>	61.7	63.5	62.5	74.4	58.4	65.4	40.7	54.6	46.7	59.5	68.8	63.8	40.1	76.1	52.5	51.6	67.8	58.6	55.2	65.0	58.3±0.5
<i>DANN</i>	67.3	51.9	58.6	80.8	46.6	59.1	59.5	49.3	53.9	73.5	55.2	63.1	57.6	56.7	57.1	66.2	45.0	53.6	66.3	49.6	56.2±0.6
<i>STA</i>	67.5	66.7	67.1	77.1	55.4	64.5	49.9	61.1	54.5	67.0	66.7	66.8	51.4	65.0	57.4	61.8	59.1	60.4	61.8	63.3	61.1±0.3
<i>OSBP</i>	66.1	67.3	66.7	76.3	68.6	72.3	48.0	63.0	54.5	72.0	69.2	70.6	59.4	70.3	64.3	60.0	62.7	64.7	64.1	66.3	64.7±0.2
<i>UAN</i>	77.5	0.1	0.2	85.0	0.1	0.1	66.2	0.0	0.0	80.6	0.1	0.2	70.5	0.0	0.0	74.0	0.1	0.2	75.2	0.0	0.1±0.0
<i>ROS</i>	67.0	70.8	68.8	72.0	80.0	75.7	51.5	73.0	60.4	65.3	72.2	68.6	53.6	65.5	58.9	59.8	71.6	65.2	61.6	72.4	66.2±0.3
<i>Ours</i>	66.8	71.8	69.2	72.5	80.1	76.1	54.5	74.2	62.9	69.4	73.3	71.3	54.7	72.1	62.2	63.7	75.3	69.0	62.3	75.8	68.0±0.4

4.3. Analysis

285 **Ablation Study.** To investigate how our method benefits known/unknown separation, we compare the performance of our Stage I with Stage I of ROS and STA. Both ROS and STA include two stages: they use a multi-rotation classifier and a multi-binary classifier to distinguish known and unknown target samples, respectively. We compute the *area under receiver operating characteristic curve* (AUC-ROC) over the normality scores \mathcal{N} on Office-31 to evaluate the performance. As shown in Table 3, the AUC-ROC of our method (93.0) is higher than 290 that of the multi-rotation used by ROS (91.5) and the multi-binary used by STA (79.9). Table 3 also reports the performance of Stage I when alternatively removing the module of deconfounding domain gaps (No DDP), the module of 295 ensembling multiple transformations (No EMT), and the label smoothing (No LS). The performance of all above cases drops significantly compared to our complete method, verifying each component’s importance: (1) the DDP module can shield the distractions caused by confounding styles and contexts from source domain during separation; (2) the EMT module can produce reliable 300 normality scores by the calibration from the ensemble; (3) label smoothing is helpful to suppress the overconfident predictions. To verify the efficiency of the self-ensembling method in OSDA, we also compare our method with the wildly adopted GRL [16] based on our Stage I. Table 3 shows that self-ensembling

Table 3: Ablation analysis.

STAGE I (AUC-ROC)	A \rightarrow W	A \rightarrow D	D \rightarrow W	W \rightarrow D	D \rightarrow A	W \rightarrow A	Avg.
Ours	91.2	91.1	99.6	99.7	89.9	86.2	93.0
Multi-Rotation (from ROS)	90.1	88.1	99.4	99.9	87.5	83.8	91.5
Multi-Binary (from STA)	83.2	84.1	86.8	72.0	75.7	78.3	79.9
Ours - No DDP	84.6	83.9	90.8	80.4	81.3	83.5	84.1
Ours - No EMT	88.4	87.9	99.0	99.6	84.7	83.9	90.6
Ours - No LS	89.8	89.1	98.4	99.7	87.5	86.9	91.9
STAGE II (HOS)	A \rightarrow W	A \rightarrow D	D \rightarrow W	W \rightarrow D	D \rightarrow A	W \rightarrow A	Avg.
Ours	83.9	84.5	98.0	99.4	79.5	78.9	87.4
Ours Stage I - GRL Stage II	84.6	84.0	98.4	99.3	79.1	76.4	87.0
Ours Stage I - No DDP in Stage II	83.3	83.9	98.2	99.4	78.6	77.8	86.9

outperforms GRL by 0.4% in average. Furthermore, we also evaluate the role of the DDP in Stage II. As shown in Table 3, our full method outperforms the case when removing the DDP module (No DDP in stage II), which verifies the proposed DDP is also helpful for domain adaptation.

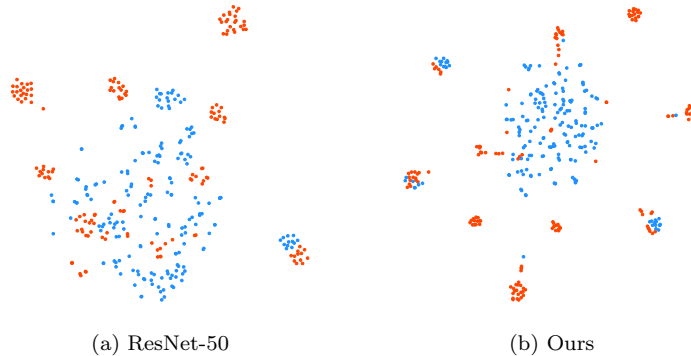


Figure 4: Visualization of obtained target features for $A \rightarrow D$ using t-SNE. Known and unknown samples are denoted as red and blue points respectively. (a): Feature obtained by ResNet-50. (b): Feature obtained by our method. Best viewed in color.

Feature Visualization. To intuitively showcase the effectiveness of our method, we visualize features of target samples from the ResNet-50 and our method on the $A \rightarrow D$ task by t-SNE [51]. The features obtained by ResNet-50 can be served as the initial state without adaptation. As shown in Fig.

4a, the features of unknown classes and several known classes mix together, demonstrating that ResNet-50 cannot separate known and unknown classes. In Figure 4b, our method is capable of separating known and unknown features and discriminating different known classes.

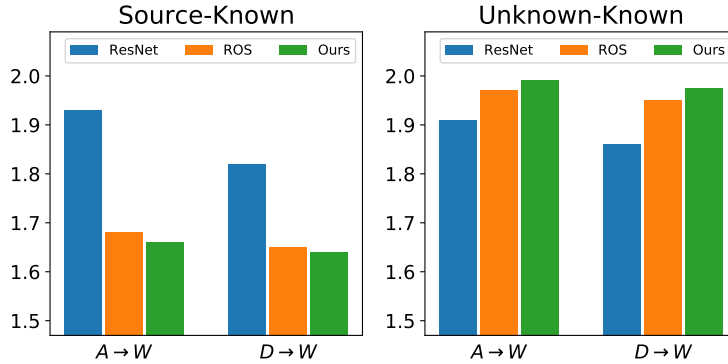


Figure 5: The value of \mathcal{A} -distance for source-known (smaller is better) and unknown-known (larger is better).

Distribution Discrepancy. As discussed in [52], distribution discrepancy can be measured by the \mathcal{A} -distance. It is defined as $d_{\mathcal{A}} = 2(1 - 2\epsilon)$, where ϵ indicates the generalization error of a domain classifier. A larger distribution discrepancy corresponds with a larger $d_{\mathcal{A}}$ and vice versa. We compute $d_{\mathcal{A}}$ for both *source-known* and *known-unknown*: *source-known* represents the distribution discrepancy between source samples and target known samples, and *known-unknown* represents the distribution discrepancy between target known samples and target unknown samples. We compare our method with ResNet-50 and ROS using a kernel SVM as the classifier on two tasks $A \rightarrow W$ and $D \rightarrow W$. From the Fig. 5, we can observe that $d_{\mathcal{A}}$ for *source-known* using our method is much smaller than the ResNet-50 (source-only), while that for *known-unknown* is larger than ResNet-50. The above observations demonstrate that our method can align the source and target known data while filtering out target unknown samples.

Sensitivity to Varying Openness. The openness is defined as $\mathbb{O} = 1 - \frac{|C_s|}{|C_t|}$, and the value of openness in the standard OSDA setting is around 0.5

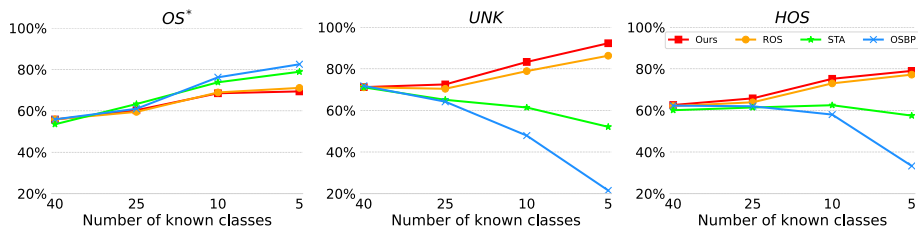


Figure 6: Accuracy (%) over the four different openness levels.

which means the number of the known and unknown target classes is close. For example, the openness of Office-31 is $\mathbb{O} = 1 - \frac{10}{21} = 0.52$ and that of Office-Home is $\mathbb{O} = 1 - \frac{25}{65} = 0.62$. In practical applications, the number of unknown target classes may exceed the number of known classes by a large margin, with openness approaching 1. To testify the robustness of our method, we conduct experiments on Office-Home with the following different openness levels: $\mathbb{O} = 0.38$ (40 known classes), $\mathbb{O} = 0.62$ (25 known classes), $\mathbb{O} = 0.85$ (10 known classes), $\mathbb{O} = 0.92$ (5 known classes). As shown in Fig. 6, the performance of OSBP and STA drops a lot with larger \mathbb{O} , as they are unable to reject the unknown instances well. In contrast, our method and ROS are resistant to the change in openness. In addition, our method outperforms ROS consistently, owing to its advanced ability of separating between known and unknown samples.

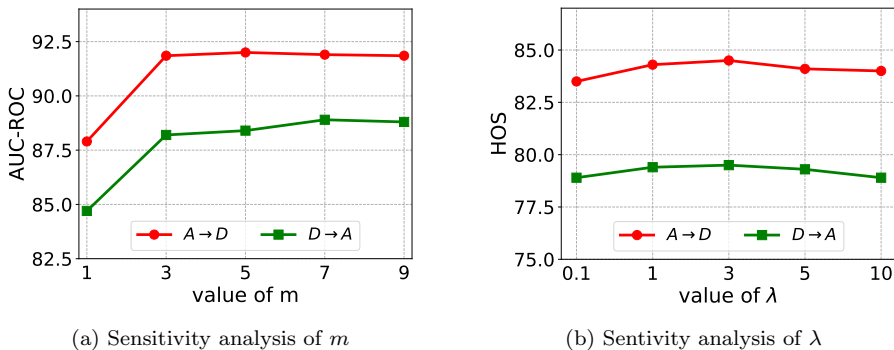


Figure 7: Hyper-parameters sensitivity analysis.

Sensitivity to Hyper-parameters. We investigate the sensitivity of two
345 hyper-parameters: the number of transformations m (in Eq. 4 and 5) and the
trade-off weight λ (in Eq. 12). The experiments are performed on two tasks
 $A \rightarrow D$ and $D \rightarrow A$ with ResNet-50 as the backbone. We plot the relationship
of the $AUC - ROC$ and the value of m in Fig. 7a, and the relationship of the
 HOS and the value of λ in Fig. 7b. Specifically, $m = 0$ denotes the ablation
350 where EMT is not used. We can observe that our method is not sensitive to
both hyper-parameters. We underline that the same hyper-parameters are used
for all 18 domain pairs demonstrating that the choice of the hyperparameters’
value is robust across datasets.

5. Conclusions

355 In this paper, we propose a novel OSDA method that can conduct effective
known and unknown separation. Specifically, we propose an explicit module of
deconfounding domain gaps (DDP) that enables the model to recognize a class
as known (or unknown) because of the class semantics rather than the confusion
caused by spurious styles or contexts. In addition, to obtain the reliable nor-
360 mality scores, we also propose a module of ensembling multiple transformations
(EMT) to calibrate the model output. The accurate known/unknown separa-
tion results boost the overall performance of the OSDA model. Experimental
results on two standard datasets show that the proposed method outperforms
the state-of-the-art OSDA methods, especially with a large margin on recogniz-
365 ing unknown samples.

Acknowledgments

This research is supported by the Agency for Science, Technology and Re-
search (A*STAR) under its AME YIRG Grant (Project No. A20E6c0101),
the National Key Research and Development Program of China (No. 2020
370 YFA0714103), the Science & Technology Development Project of Jilin Province,
China (20190302117GX), the Innovation Capacity Construction Project of Jilin

Province Development and Reform Commission (2019C053-3), and Graduate Innovation Fund of Jilin University (101832020CX179).

References

- 375 [1] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- [2] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs, IEEE transactions on pattern analysis and
380 machine intelligence 40 (4) (2017) 834–848.
- [3] K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask r-cnn, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2961–2969.
- 385 [4] A. Gretton, A. Smola, J. Huang, M. Schmittfull, K. Borgwardt, B. Schölkopf, Covariate shift and local learning by distribution matching, MIT Press, Cambridge, MA, USA, 2009, pp. 131–160.
- [5] S. J. Pan, Q. Yang, A survey on transfer learning, IEEE Transactions on knowledge and data engineering 22 (10) (2010) 1345–1359.
- 390 [6] P. Panareda Busto, J. Gall, Open set domain adaptation, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 754–763.
- [7] K. Saito, S. Yamamoto, Y. Ushiku, T. Harada, Open set domain adaptation by backpropagation, in: Proceedings of the European Conference on
395 Computer Vision (ECCV), 2018, pp. 153–168.
- [8] H. Liu, Z. Cao, M. Long, J. Wang, Q. Yang, Separate to adapt: Open set domain adaptation via progressive separation, in: Proceedings of the

IEEE/CVF Conference on Computer Vision and Pattern Recognition,
2019, pp. 2927–2936.

- 400 [9] T. Wang, C. Zhou, Q. Sun, H. Zhang, Causal attention for unbiased visual
recognition, in: International Conference on Computer Vision, 2021.
- [10] G. French, M. Mackiewicz, M. Fisher, Self-ensembling for visual domain
adaptation, in: International Conference on Learning Representations,
2018.
- 405 [11] M. Long, Y. Cao, J. Wang, M. Jordan, Learning transferable features with
deep adaptation networks, in: International Conference on Machine Learning,
2015, pp. 97–105.
- [12] M. Long, H. Zhu, J. Wang, M. I. Jordan, Unsupervised domain adaptation
with residual transfer networks, in: Proceedings of the 30th International
410 Conference on Neural Information Processing Systems, 2016, pp. 136–144.
- [13] M. Long, H. Zhu, J. Wang, M. I. Jordan, Deep transfer learning with joint
adaptation networks, in: Proceedings of the 34th International Conference
on Machine Learning-Volume 70, JMLR. org, 2017, pp. 2208–2217.
- [14] C. Chen, Z. Fu, Z. Chen, S. Jin, Z. Cheng, X. Jin, X.-S. Hua, Homm:
415 Higher-order moment matching for unsupervised domain adaptation, in:
Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34,
2020, pp. 3422–3429.
- [15] N. Courty, R. Flamary, D. Tuia, A. Rakotomamonjy, Optimal transport
for domain adaptation, IEEE transactions on pattern analysis and machine
420 intelligence 39 (9) (2016) 1853–1865.
- [16] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Lavio-
lette, M. Marchand, V. Lempitsky, Domain-adversarial training of neural
networks, The Journal of Machine Learning Research 17 (1) (2016) 2096–
2030.

- 425 [17] E. Tzeng, J. Hoffman, K. Saenko, T. Darrell, Adversarial discriminative domain adaptation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 7167–7176.
- [18] M. Long, Z. Cao, J. Wang, M. I. Jordan, Conditional adversarial domain adaptation, in: Advances in Neural Information Processing Systems, 2018, 430 pp. 1640–1650.
- [19] S. Xie, Z. Zheng, L. Chen, C. Chen, Learning semantic representations for unsupervised domain adaptation, in: International Conference on Machine Learning, 2018, pp. 5419–5428.
- [20] K. Saito, K. Watanabe, Y. Ushiku, T. Harada, Maximum classifier discrepancy for unsupervised domain adaptation, in: Proceedings of the IEEE 435 Conference on Computer Vision and Pattern Recognition, 2018, pp. 3723–3732.
- [21] C. Chen, W. Xie, W. Huang, Y. Rong, X. Ding, Y. Huang, T. Xu, J. Huang, Progressive feature alignment for unsupervised domain adaptation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 440 2019, pp. 627–636.
- [22] H. Tang, K. Jia, Discriminative adversarial domain adaptation, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34, 2020, pp. 5940–5947.
- 445 [23] S. Cui, S. Wang, J. Zhuo, C. Su, Q. Huang, Q. Tian, Gradually vanishing bridge for adversarial domain adaptation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 12455–12464.
- [24] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, D. Krishnan, Unsupervised pixel-level domain adaptation with generative adversarial networks, 450 in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 3722–3731.

- [25] Y. Taigman, A. Polyak, L. Wolf, Unsupervised cross-domain image generation, in: ICLR, 2016.
- 455 [26] M.-Y. Liu, T. Breuel, J. Kautz, Unsupervised image-to-image translation networks, in: Advances in neural information processing systems, 2017, pp. 700–708.
- [27] G. Yang, H. Xia, M. Ding, Z. Ding, Bi-directional generation for unsupervised domain adaptation, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34, 2020, pp. 6615–6622.
- 460 [28] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, D. Erhan, Domain separation networks, in: Proceedings of the 30th International Conference on Neural Information Processing Systems, 2016, pp. 343–351.
- [29] F. M. Carlucci, A. D’Innocente, S. Bucci, B. Caputo, T. Tommasi, Domain generalization by solving jigsaw puzzles, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 2229–2238.
- 465 [30] M. Ghifary, W. B. Kleijn, M. Zhang, D. Balduzzi, W. Li, Deep reconstruction-classification networks for unsupervised domain adaptation, in: European Conference on Computer Vision, Springer, 2016, pp. 597–613.
- 470 [31] X. Zhao, S. Wang, Adversarial learning and interpolation consistency for unsupervised domain adaptation, IEEE Access 7 (2019) 170448–170456.
- [32] Q. Feng, G. Kang, H. Fan, Y. Yang, Attract or distract: Exploit the margin of open set, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 7990–7999.
- 475 [33] S. Bucci, M. R. Loghmani, T. Tommasi, On the effectiveness of image rotation for open set domain adaptation, in: European Conference on Computer Vision, Springer, 2020, pp. 422–438.

- [34] K. You, M. Long, Z. Cao, J. Wang, M. I. Jordan, Universal domain adaptation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 2720–2729.
- [35] B. Fu, Z. Cao, M. Long, J. Wang, Learning to detect open classes for universal domain adaptation, in: European Conference on Computer Vision, Springer, 2020, pp. 567–583.
- [36] S. Zhai, Y. Cheng, W. Lu, Z. Zhang, Deep structured energy based models for anomaly detection, in: International Conference on Machine Learning, PMLR, 2016, pp. 1100–1109.
- [37] B. Zong, Q. Song, M. R. Min, W. Cheng, C. Lumezanu, D. Cho, H. Chen, Deep autoencoding gaussian mixture model for unsupervised anomaly detection, in: International conference on learning representations, 2018.
- [38] Y. Xia, X. Cao, F. Wen, G. Hua, J. Sun, Learning discriminative reconstructions for unsupervised outlier removal, in: Proceedings of the IEEE International Conference on Computer Vision, 2015, pp. 1511–1519.
- [39] C. Zhou, R. C. Paffenroth, Anomaly detection with robust deep autoencoders, in: Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining, 2017, pp. 665–674.
- [40] S. Liang, Y. Li, R. Srikant, Enhancing the reliability of out-of-distribution image detection in neural networks, in: International Conference on Learning Representations, 2018.
- [41] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, M. Kloft, Deep one-class classification, in: International conference on machine learning, PMLR, 2018, pp. 4393–4402.
- [42] K. Zhou, Y. Yang, Y. Qiao, T. Xiang, Domain generalization with mixstyle, in: International Conference on Learning Representations, 2021.
URL <https://openreview.net/forum?id=6xHJ37MVxxp>

- [43] H. Nam, H. Lee, J. Park, W. Yoon, D. Yoo, Reducing domain gap by reducing style bias, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 8690–8699.
- [44] X. Huang, S. Belongie, Arbitrary style transfer in real-time with adaptive instance normalization, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 1501–1510.
- [45] T. Karras, S. Laine, T. Aila, A style-based generator architecture for generative adversarial networks, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 4401–4410.
- [46] R. Müller, S. Kornblith, G. Hinton, When does label smoothing help?, arXiv preprint arXiv:1906.02629.
- [47] K. Saenko, B. Kulis, M. Fritz, T. Darrell, Adapting visual category models to new domains, in: European conference on computer vision, Springer, 2010, pp. 213–226.
- [48] H. Venkateswara, J. Eusebio, S. Chakraborty, S. Panchanathan, Deep hashing network for unsupervised domain adaptation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 5018–5027.
- [49] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A large-scale hierarchical image database, in: 2009 IEEE conference on computer vision and pattern recognition, Ieee, 2009, pp. 248–255.
- [50] M. L. Junguang Jiang, Bo Fu, Transfer-learning-library, <https://github.com/thuml/Transfer-Learning-Library> (2020).
- [51] L. v. d. Maaten, G. Hinton, Visualizing data using t-sne, Journal of machine learning research 9 (Nov) (2008) 2579–2605.
- [52] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, J. W. Vaughan, A theory of learning from different domains, Machine learning 79 (1) (2010) 151–175.